

Northumbria Research Link

Citation: Oktay, Kaan, Santaliz-Casiano, Ashlie, Patel, Meera, Marino, Natascia, Storniolo, Anna Maria V., Torun, Hamdi, Acar, Burak and Madak Erdogan, Zeynep (2020) A Computational Statistics Approach to Evaluate Blood Biomarkers for Breast Cancer Risk Stratification. *Hormones and Cancer*, 11 (1). pp. 17-33. ISSN 1868-8497

Published by: Springer

URL: <https://doi.org/10.1007/s12672-019-00372-3> <<https://doi.org/10.1007/s12672-019-00372-3>>

This version was downloaded from Northumbria Research Link:
<http://nrl.northumbria.ac.uk/id/eprint/41786/>

Northumbria University has developed Northumbria Research Link (NRL) to enable users to access the University's research output. Copyright © and moral rights for items on NRL are retained by the individual author(s) and/or other copyright owners. Single copies of full items can be reproduced, displayed or performed, and given to third parties in any format or medium for personal research or study, educational, or not-for-profit purposes without prior permission or charge, provided the authors, title and full bibliographic details are given, as well as a hyperlink and/or URL to the original metadata page. The content must not be changed in any way. Full items must not be sold commercially in any format or medium without formal permission of the copyright holder. The full policy is available online: <http://nrl.northumbria.ac.uk/policies.html>

This document may differ from the final, published version of the research and has been made available online in accordance with publisher policies. To read and/or cite from the published version of the research, please visit the publisher's website (a subscription may be required.)

A Computational Statistics Approach to Evaluate Blood Biomarkers for

Breast Cancer Risk Stratification

Kaan Oktay (kaan.oktay@boun.edu.tr)^{1*}, Ashlie Santaliz-Casiano
(ashlies2@illinois.edu)^{2*}, Meera Patel (patel88253@gmail.com)³, Natascia Marino
(marinon@iu.edu)^{3,4}, Anna Maria V. Storniolo (astornio@iu.edu)^{3,4}, Hamdi Torun
(hamdi.torun@northumbria.ac.uk)⁵, Burak Acar (acarbu@boun.edu.tr)¹, Zeynep Madak
Erdogan^{2,6,7,8,9,10}

¹VAVlab, Electrical & Electronics Engineering Department, Bogazici University, Istanbul, Turkey

² Division of Nutritional Sciences, University of Illinois, Urbana-Champaign, IL USA

³Susan G. Komen Tissue Bank at the IU Simon Cancer Center, Indianapolis, IN

⁴ Department of Medicine, Indiana University School of Medicine, Indianapolis, IN

⁵Faculty of Engineering and Environment, University of Northumbria, Newcastle upon Tyne, UK

⁶ Department of Food Sciences and Human Nutrition, University of Illinois, Urbana-Champaign, IL USA

⁷ National Center for Supercomputing Applications, University of Illinois, Urbana-Champaign, Urbana, IL, USA

⁸ Cancer Center at Illinois, University of Illinois, Urbana-Champaign, Urbana, IL, USA

⁹ Beckman Institute for Advanced Science and Technology, University of Illinois at Urbana-Champaign, Urbana, IL, USA

¹⁰ Carl R. Woese Institute for Genomic Biology, University of Illinois, Urbana-Champaign, Urbana, IL, USA

**These authors contributed equally*

Corresponding Author and Lead Contact: Zeynep Madak-Erdogan,

zmadake2@illinois.edu

Abstract

Breast cancer is the second leading cause of cancer mortality among women. Mammography and tumor biopsy followed by histopathological analysis are the current methods to diagnose breast cancer. Mammography does not detect all breast tumor subtypes; especially those arise in younger women or women with dense breast tissue, and are more aggressive. There is an urgent need to find circulating prognostic molecules and liquid biopsy methods for breast cancer diagnosis and reducing the mortality rate. In this study, we systematically evaluated metabolites and proteins in blood to develop a pipeline to identify potential circulating biomarkers for breast cancer risk. Our aim is to identify a group of molecules to be used in the design of portable and low-cost biomarker detection devices. We obtained plasma samples from women who are cancer free (healthy) and women who were cancer free at the time of blood collection but developed breast cancer later (susceptible). We extracted potential prognostic biomarkers for breast cancer risk from plasma metabolomics and proteomics data using statistical and discriminative power analyses. We pre-processed the data to ensure the quality of subsequent analyses, and used two main feature selection methods to determine the importance of each molecule. After further feature elimination based on pairwise dependencies, we measured the performance of logistic regression classifier on the remaining molecules and compared their biological relevance. We identified six signatures that predicted breast cancer risk with different specificity and selectivity. The best performing signature had 13 factors. We validated the difference in level of one of biomarkers, SCF/KITLG, in plasma from healthy and susceptible individuals. These biomarkers will be used to develop low-cost liquid biopsy methods towards early

1
2
3
4 identification of breast cancer risk and hence decreased mortality. Our findings provide
5
6 the knowledge basis needed to proceed in this direction.
7
8

9
10 **Keywords:** Liquid biopsy, Breast cancer risk, circulating biomarker, Machine
11
12 learning, Feature selection
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Background

Breast cancer is the second leading cause of death among adult women.

According to World Health Organization, there is a sharp rise in overall number of breast cancer incidences world-wide due to changes in life style, reproductive factors and increased life expectancy [1]. Fifty eight percent of all breast cancer-related deaths occur in middle- and low-income countries. While survival rates for breast cancer are around 80% in developed countries, this rate decreases to 60% in middle-income and to 40% in low-income countries due to lack of early detection programs leading to diagnoses in late stages, where 80% of these tumors are incurable [2, 3]. In the middle- and low-income countries, mammography and other expensive and technologically complicated methods are unattainable due to high costs and shortage of trained personnel. [4, 5] Moreover, mammograms are more likely to detect ER-positive breast cancer [6] and are not recommended for younger women. In addition, diagnosis at an earlier stage using conventional procedures is not prognostic for all race groups, for example, the probability of an African-American woman with small-sized tumors presenting with metastasis is higher than that of a Caucasian women.[7] Thus, there is a critical need for affordable, portable and accurate means of detecting breast cancer risk before the tumors arise. Development of such technologies has the potential to expedite the solution for the growing health problem to prevent increasing death and disability among women especially in low- and middle-income countries.

Currently, a handful of biomarkers are used in the clinic for breast cancer diagnosis. These biomarkers are proteins overexpressed in certain subtypes of breast tumors and help clinicians plan treatment. Up to date, a limited number of breast

cancer biomarkers demonstrated clinical utility, including Estrogen Receptor alpha (ER α), Progesterone Receptor (PgR), [8] and human epidermal growth factor receptor 2 (HER2) to predict effectiveness of systemic therapy and the Oncotype DX-21 gene score to predict benefits of chemotherapy. [9-11] Studies evaluating other predictive biomarkers are in progress for Breast Cancer susceptibility genes (BRCA1 and BRCA2) circulating tumor cells (CTCs), HER2 (+), TOP2A (in subjects with HER2 overexpression) and HER2 (when is negative in tumors but is positive in the CTCs). [12] Circulating tumor DNA (ctDNA) is increasingly used in the clinic, particularly for advanced solid tumors. [13-15] However, clinical utility and validity of ctDNA assays in early stage cancers is not as clear. [15] Further, we still lack reliable biomarkers to detect breast cancer risk before the tumors arise. Lack of such biomarkers hinders establishment of reliable screening or prevention programs.

To address this critical need, we systematically evaluated metabolites and proteins in plasma to identify potential biomarkers for breast cancer risk that can be utilized to develop minimally invasive, affordable, portable, and accurate screening devices. In this study, our focus is on liquid biopsy samples from plasma that have the potential to provide simple and minimally invasive information for diagnostic decisions. We developed an efficient pipeline to analyze liquid biopsy samples, to detect blood biomarkers and to identify the risk for breast cancer before tumors arise. This pipeline paves the way towards developing the aforementioned screening devices to be used in the field by basic level healthcare workers in low-resource environments.

Methods

Patients and Plasma Samples

1
2
3
4 All studies were approved by the Indiana University Institutional Review Board
5
6 (IRB protocol number 1011003097). All research was carried out in compliance with the
7
8 Helsinki Declaration. Donors provided broad written consent for the use of their
9
10 specimens in research. The written consent document informed the donor that the
11
12 donated specimens and medical data would be used for the general purpose of helping
13
14 to determine how breast cancer develops. It was explained in the written consent that
15
16 the exact laboratory experiments were unknown at the time of donation, and that
17
18 proposals for use of the specimens would be reviewed and approved by a panel of
19
20 independent researchers before specimens and/or data were released for research
21
22 purposes. Hematoxylin and eosin stained sections of the FFPE tissue of the identified
23
24 donors were reviewed by pathologist to confirm the absence of histological
25
26 abnormalities. In order to exclude or control confounding variables such as age, racial
27
28 and ethnic background and menopausal status the subjects in the two cohorts,
29
30 susceptible and healthy controls, were matched by selection of the comparison group
31
32 (healthy controls) with respect to the distribution of the above mentioned confounders in
33
34 susceptible group.
35
36
37
38
39
40
41
42

43 ***Plasma preparation***

44

45
46 Blood was drawn into the Plasma Separator tube (Vacutainer Venous Blood
47
48 Collection Tubes: SST* Plasma Separation Tube, Fisher Scientific cat. #0268396) and
49
50 gently mixed by inverting the tube 5 times. Forty-five minutes (± 10 min.) after the blood
51
52 had been drawn, the Plasma Separator Tube was placed into a minicentrifuge
53
54 (Eppendorf centrifuge 5702) and centrifuged at 1200 rcf for ten minutes at room
55
56
57
58
59
60
61
62
63
64
65

temperature. A repeater pipet was used to aliquot 600ul of the plasma into each of five cryogenic vials. Samples were stored at -80°C until use.

OLINK Protein biomarker and whole metabolite profiling assays

All the samples from human studies were handled and analyzed in accordance with UIUC IRB protocol #06741 and as previously described [16]. 10 µl of plasma samples from Komen Tissue Bank were submitted to OLINK biosciences for cancer and inflammation biomarker analysis. 50 µl of plasma samples were submitted to the Metabolomics Center at UIUC. GC/MS whole metabolite profiling was performed to detect and quantify the metabolites by using gas chromatography-mass spectrometry (GC/MS) analysis. Metabolites were extracted from 50 µl of plasma according to Agilent Inc. application notes. The hentriacontanoic acid was added to each sample as the internal standard prior to derivatization. Metabolite profiles were acquired using an Agilent GC/MS system (Agilent 7890 gas chromatograph, an Agilent 5975 MSD, and an HP 7683B autosampler). The spectra of all chromatogram peaks were evaluated using the AMDIS 2.71 and a custom-built database with 460 unique metabolites. All known artificial peaks were identified and removed prior to data analysis. To allow the comparison between samples, all data were normalized to the internal standard in each chromatogram.

Statistical Analysis

Preprocessing of Measurements

We normalized all individuals' plasma data in each dataset with respect to the healthy individuals' data in the respective dataset to factor out potential differences in

1
2
3
4 data acquisition. More specifically, we performed the following procedure separately for
5
6 both datasets. For each molecule in a dataset, we subtracted the mean measurement of
7
8 that molecule in healthy individuals from all individuals' measurements and divided this
9
10 difference by the standard deviation of that molecule's measurements in healthy
11
12 individuals. Thus, we converted each single measurement to a z-score which describes
13
14 the deviation of that measurement from the mean of healthy individual's, in terms of the
15
16 standard deviation among healthy individuals. As the final step, we merged two
17
18 datasets, which were normalized with respect to their own healthy individuals, and
19
20 obtained a dataset with 49 susceptible and 47 healthy individuals.
21
22
23
24
25

26 ***Molecule Ranking, Elimination and Performance Assessment***

27
28
29 A two-stage procedure is applied to identify the molecule sets with high
30
31 discriminative power between the healthy and the susceptible groups. The first stage
32
33 involves ranking all molecules with respect to their individual discriminative powers
34
35 (importance ranking). The second stage involves molecule elimination (selection) based
36
37 on their interdependencies.
38
39
40
41

42 To independently assess each of 181 molecules, we used two different methods.
43
44 In the first method, we applied Student's t-test to test the null hypothesis that the
45
46 measurements in the two groups come from the same distribution. All molecules were
47
48 ranked based on the corresponding p-values to get a short-list of the top-ranking 20
49
50 molecules with the lowest p-values, discarding the others from further processing. In the
51
52 second method, we applied the random forest algorithm to assess the discriminative
53
54 power of each of the 181 molecules individually by using the mean decrease impurity
55
56 (Gini importance), which is defined as the mean decrease in node impurity over all the
57
58
59
60
61
62
63
64
65

1
2
3
4 trees in the forest. This time, all molecules were ranked based on their Gini importance
5 values to get the top-ranking 20 molecules with the highest importance values. No further
6
7 threshold was applied to these top-ranking molecules at this stage for both methods, as
8
9 the low-ranking molecules in these lists may potentially have significant marginal
10
11 contribution to a subset of molecules when used together.
12
13
14
15

16
17 To generate an optimum subset of the top 20 molecules identified by Student's t-
18
19 test or random forest, we used the following iterative procedure. We initialized a "selected
20
21 molecules" list (S-list) with the top-ranking molecule and an "unselected molecules" list
22
23 (U-list) with the remaining 19 ranked molecules. We iteratively assessed the individual
24
25 molecules in the U-list with respect to the molecules set represented by the S-list, and
26
27 added the ones that have a positive contribution to the S-list while discarding the others.
28
29 Three different approaches are applied to assess whether a molecule has a positive
30
31 contribution to the S-list: (i) Manual selection: Logistic Regression (LR) classifiers, to
32
33 identify healthy and susceptible groups, are trained and tested iteratively by using the *selected*
34
35 molecules (S-list) and the top-ranking *unselected* molecule (U-list) as the features. The classifier
36
37 performance is assessed using the *selected* molecules' AUC (Area Under Curve) of ROC
38
39 (Receiver Operator Characteristic) curves. After each iteration, if the AUC is increased, the top-
40
41 ranking *unselected* molecule is added to the S-list, otherwise discarded. The iterations stop
42
43 when the U-list is exhausted. (ii) Paired t-test: The inter-molecule dependencies, as
44
45 measured by the paired t-test, is used to select the molecules from the U-list to be added
46
47 to the S-list. We first computed the paired t-test p-values for each pair of molecules among
48
49 the aforementioned top-ranking 20 molecules with the null hypothesis being that both
50
51 come from the same distribution. Using these p-values, we iteratively discarded the
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 molecules from the U-list that have p-value larger than 0.05 when tested with any of
5
6 the molecules from the S-list and moved the *unselected* molecule from U-list to S-list with
7
8 the lowest maximum p-value (<0.05) when tested with the *selected* molecules. The
9
10 iterations stop when the U-list is exhausted. (iii) Correlation Analysis: The second approach
11
12 described above is repeated by replacing the null hypothesis testing with the correlation analysis
13
14 as measured by the Pearson's correlation coefficient (pCC). We used 0.5 as the pCC
15
16 threshold.
17
18
19
20

21
22 Finally, we performed LR classification (4-fold cross-validation with 500 iterations)
23
24 using the top-ranking N molecules in each list, where N runs from one to the length of the
25
26 corresponding list. Of note, use of LR for performance assessment of classification at this
27
28 last step is distinct from the earlier use of LR for manual selection of the molecules.
29
30
31

32 **SCF/KITLG quantification using enzyme-linked immunosorbent assay**

33 **(ELISA)**

34
35
36

37 Plasma samples from both groups were collected, and stored at -80°C until the
38
39 time of assay. We used an enzyme-linked immunosorbent assay (ELISA) kit for
40
41 SCF/KITLG (Sigma, catalogue no. RAB0330). Samples were diluted two fold per
42
43 suggestion from the manufacturer. For SCF/KITLG the antibody, concentrate was
44
45 diluted 100 fold with 1X Diluent Buffer. To prepare the HRP-Streptavidin Concentrate
46
47 the vial was spin and diluted 400 times with 1X Diluent Buffer. A 50ng/ml stock solution
48
49 was used to make the standard curve: 2000 pg/ml, 666.7pg/ml, 222.2 pg/ml,
50
51 74.07pg/ml, 24.69pg/ml, 8.23 pg/ml and 2.74pg/ml for SCF/KITLG. The Human
52
53 SCF/KITLG antibody-pre coated ELISA wells were filled with 100 μl of either serially
54
55 diluted standard protein and plasma samples. After 2.5 hours incubation with gentle
56
57
58
59
60
61
62
63
64
65

shaking at room temperature, 100µl of 1x SCF/KITLG Biotinylated Detection Antibody were added to the wells. After one-hour incubation with shaking at room temperature, the solution was discarded and the wells were washed four times using 300µL wash buffer solution. Final wash was aspirated and plates were inverted to remove any remaining was buffer. Then, 100µl of Prepared HRP Streptavidin solution was added to each well, and incubated for 45 minutes at room temperature with gentle shaking. The solution was discarded and washed four times as described previously. 100µl of ELISA colorimetric TMB reagent was added to each well and incubated for 30 minutes at room temperature in the dark. After this, 50µl of Stop solution was added to each well. Immediately after color development, the OD values were measured at 450nm using Cytation 5 Cell Imaging Multi Mode- Reader (Biotek) and SCF/KITLG concentrations were calculated from specific calibration curves prepared with known standard solutions. Diluent buffer served as blank and the OD of these wells was subtracted from the values.

Results

Identification of circulating factor signatures for future breast cancer risk assessment

Because we wanted to identify circulating factors that might indicate future breast cancer risk, we utilized plasma samples from a cohort of healthy controls (Healthy) and individuals who were clinically healthy at the time of plasma collection but later had a diagnosis of breast cancer (Susceptible). We analyzed plasma samples using whole metabolite profiling and OLINK biomarker analysis for a panel of inflammation and cancer-related proteins. We used two different sample sets, one with 39 susceptible and

36 healthy and the other with 10 susceptible and 11 healthy individuals, which were collected at different times. In the first set, 22 out of 39 susceptible and 23 out of 36 healthy individuals were postmenopausal status and remaining ones were premenopausal. In the second dataset, 7 out of 10 susceptible and 8 out of 11 healthy individuals were postmenopausal status and remaining ones were premenopausal. Average time to diagnosis was 3.7 years after samples donation (median is 4 years). Data from two datasets were pre-processed separately because they were acquired at different times, and were expected to have a variation due to external factors. Plasma levels of 295 different molecules for the first dataset and 339 different molecules for the second dataset were detected for the individuals. Some molecules had missing values (were not detected by metabolomics or OLINK approach) for some individuals and further, some molecules were not measured for both datasets. All these molecules were excluded from the analysis. Therefore, we analyzed 181 different molecules, consisting of metabolites and proteins, which have plasma level values for every subject in both datasets.

In order to generate an inclusive list of features that would best discriminate between healthy and susceptible individuals, we took a stepwise approach where we first screened all molecules that contribute to increased classifier performance (LR) and then iteratively eliminate the redundant ones for both top-ranking molecule lists obtained by either of the initial molecule selection methods. We initially selected two different groups of 20 molecules (out of 181 molecules) using the Student's t-test and random forest (600 trees) methods to rank all 181 molecules with respect to their (healthy vs susceptible) discriminative power. Top-ranking 20 molecule sets obtained by two

different feature selection methods, Student's t-test and random forest, contain 10 common molecules, highly concentrated in the upper halves of the lists. For example, 4 out of 5 top ranking molecules are common in both datasets. To assess the pairwise dependencies among the most discriminative 20 molecules and further reduce the number of features in our lists we used the paired t-test (**Table 1-Student's t-test-, Table 2-random forest-**) or pairwise correlation analysis (**Table 3-Student's t-test-, Table 4-random forest-**). To ensure that all molecules that might positively contribute to classifier performance are included in the signature we performed logistic regression. Finally, to eliminate redundant molecules, we utilized paired t-test p-values ($p > 0.05$) and/or correlation coefficients ($pCC > 0.5$) to discard one of the molecules in that pair. Our approach resulted in six molecule signatures (**Table 5-Student's t-test-, Table 6-random forest-**).

Assessment of classification performances of molecule signatures using machine-learning approach

In order to test the classification performance of each molecule signature, we performed LR classification using the molecules indicated in Table 5 and Table 6. Of note, use of LR for performance assessment of classification at the last step is distinct from the earlier use of LR for manual selection of the molecules. Our top 20 feature list generated by Student's t-test contained MMP-10, MCP-3, SCF/KITLG, TRAIL, ENRAGE, MAD HOMOLOG 5 (SMAD5), CXL17, HK11, FGF-BP1, XPNPEP2, C15:0 (Pentadecanoic acid), PPY, FGF-5, FGF-21, ESM-1, FASLG, CD160, TNFB, CTSV and ADA (**Figure 1A**). Unsupervised clustering of the data using this list of molecules separated healthy, and susceptible individuals, only two healthy individuals were

classified with susceptible individuals and only one individual was classified together with healthy individuals (**Figure 1A**). This list without any further feature elimination achieved AUC value of 0.83 (**Figure 1B**). Reduction of feature number to 13 using manual selection increased AUC value to 0.85 ± 0.04 (**Figure 1C**). Further reduction of feature using correlation analysis (**Figure 1D**; $AUC = 0.78 \pm 0.04$) or paired t-test (**Figure 1E**; $AUC = 0.69 \pm 0.03$). On the other hand, AUC values achieved by molecule signatures using Random Forest had lower performance (**Figure 2**). This list contained XPNPEP2, phosphoric acid, FGF-BP1, MAD HOMOLOG 5, ESM-1, SCF/KITLG, TRAIL, PD-L1, FLT3L, 4E-BP1, MCP-1, PPY, FGF-5, FASLG, MMP-10, EPHA2, CD27, CXCL1, HK14 and TLR3 (**Figure 2A**). Unsupervised clustering of the data using this list of molecules was less successful in separating healthy and susceptible individuals, 10 of susceptible individuals were classified with healthy individuals (**Figure 2A**). Using all 20 factors achieved AUC of 0.80 ± 0.05 (**Figure 2B**). Reducing the molecule number to 10 using manual selection (**Figure 2C**, $AUC = 0.80 \pm 0.04$), to 11 using correlation analysis (**Figure 2D**, $AUC = 0.76 \pm 0.05$) or to two using paired t-test (**Figure 2E**, $AUC = 0.67 \pm 0.04$) did not improve the AUC values. To sum, initial feature selection using Student's t-test followed by manual selection using LR gave us the best performing list of 13 circulating molecules from plasma for differentiating between healthy and susceptible individuals.

Biological Relevance of Biomarkers

Our best performing list contained SCF/KITLG, MMP-10, MAD HOMOLOG5, CXL17, MCP-3, FGF05, FASLG, CD160, TNFB, ESM-1, FGF-21, XPNPEP2 and CTSV (**Figure 3A**). In order to increase our understanding of molecules in the best performing molecule list. Unsupervised clustering of the data using this list of molecules separated

1
2
3
4 healthy and susceptible individuals accurately (**Figure 3A**). To delve further into
5
6 direction of change in the plasma levels of identified molecules we compared level of
7
8 individual molecules in healthy vs. susceptible individuals. Six of the 13 molecules,
9
10 including SCF/KITLG, MAD HOMOLOG 5, FASLG, MMP-10, XPNPEP2 and CXL17
11
12 were statistically significantly different between the two groups (**Figure 3B**). Since our
13
14 aim was to identify the molecules that have marginal but significant contribution to the
15
16 classification task when used together with other molecules, even if they have weak
17
18 discriminative power on their own, we still included these molecules with poor t-test
19
20 performance individually, $p\text{-value} > 0.05$, or low random forest importance in the final
21
22 lists. We were particularly interested in SCF/KITLG as this molecule was the top
23
24 molecule identified in both feature selection methods (**Table 5 and 6**). Overall,
25
26 SCF/KITLG levels were lower in individuals with increased breast cancer risk (**Figure**
27
28 **3C**). We also validated our finding from OLINK analysis using another independent
29
30 method, ELISA analysis, and verified that level of this protein is lower in susceptible
31
32 individuals (**Figure 3B**).
33
34
35
36
37
38
39

40 Discussion

41
42 In this study, we developed a pipeline to identify plasma biomarkers of breast
43
44 cancer risk using a combination of classical statistics methods and machine learning
45
46 approaches, and independently validated one of the identified biomarkers, SCF/KITLG.
47
48 By iterative feature selection, elimination and performance testing we generated a
49
50 molecular signature of plasma biomarkers that can discriminate between healthy and
51
52 breast cancer susceptible individuals. Because of our approach, some of the molecules
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 in this signature had weak discriminative power on their own, yet they contributed
5
6 significantly to the discriminative power of the signature.
7
8

9 A biomarker is a biomolecule, such as DNA, RNA, proteins, hormones and
10
11 chemical modifications that can be measured to describe that an abnormal or a normal
12
13 process is taking place within the organism. [12] A cancer biomarker can arise due to
14
15 changes in the DNA (mutations), rearrangements, deletions (missing copies), or
16
17 amplifications. Biomarkers might affect various hallmarks of cancer including cell cycle,
18
19 cell death, or immunological properties of the tumor and indicate the risk of developing
20
21 cancer, its progression and response to therapy. [8, 9, 12]
22
23
24
25

26 Previously, Kazarian et al, studied pre-diagnostic samples from the UK
27
28 Collaborative Trial of Ovarian Cancer Screening (UKCTOCS). Serum samples were
29
30 taken from 239 women who were diagnosed with invasive ductal carcinoma in breast
31
32 (IDC), months to years after sample donation [17]. These patients were post-
33
34 menopausal women with ages ranging from 50 to 74, who were healthy cases at the
35
36 moment of recruitment but that later developed breast cancer. Hence, this group studied
37
38 the ability of several serum markers to detect breast cancer cases before these patients
39
40 were diagnosed. They studied CA 15-3, RANTES/CCL5, OPN, PAI-1, SLP1, HSP90A,
41
42 IGFBP3, APOC1 and PAPPA. They concluded that only 3 out of the 9 serum markers,
43
44 (CA 15-3, PAI1 and HSP90A) were potential prognostic biomarkers[17]. Those
45
46 analyses were performed using a limited panel of proteins. However, in our analysis, we
47
48 characterized more than 300 proteins and metabolites in plasma and used a final list of
49
50 181 molecules to generate our signatures.
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

One potential marker of interest we identified is SCF/KITLG protein. KITLG protein is expressed in 53% of breast cancer cell lines[18]. SCF/KITLG was shown to have a proliferative role in BCK4 cells, and when it is reduced, it decreased estrogen-induced proliferation[19]. We identified lower level of this biomarker in plasma from women with breast cancer risk. Since at the time of the blood draw the women did not have tumors it is not possible for us to infer level of this protein in the tumor. Whether SCF plays a role in the induction of breast tumors or lower plasma levels of this protein contributes to the tumor biology needs to be determined.

Several of the molecules in our signature were also implicated in cancer biology. For example, MAD HOMOLOG5/SMAD5 plays a role in breast cancer cell stemness, and resistance to chemotherapy. [20, 21] FGF-5, FASLG, CTSV and ESM-1 expression is associated with lower survival and worst outcomes. [22-27] MMP-10 affects angiogenesis and apoptosis [28, 29]; XPNPEP2 is overexpressed in cervical cancer patients and increases motility and invasiveness of tumors. [30] FGF-21, TNFB contributes to metastatic potential of breast cancer cells. [31, 32] CXL17 [33], MCP-3[34] and CD160 [35] play a role in recruitment of immune cells. TNFB/LTA polymorphisms increased the cancer risk in various populations. [36, 37] All these studies focused on the tumors or patients that already have cancers. The impact of proteins in our signature on breast cancer risk and initiation remains to be established. Direction of differences in the plasma levels of these proteins between healthy and susceptible individuals might be different from what is reported in already established tumors and might indicate a different role for these proteins at early stages of tumor development.

1
2
3
4 More recently, liquid biopsy methods supported with machine-learning
5
6 approaches have been used for the detection of different cancer types.[15, 38, 39] For
7
8 example, Cohen et al. recently demonstrated the capability of detecting eight different
9
10 cancer types including breast cancer using circulating tumor DNA (ctDNA) and protein
11
12 biomarkers [40]. They reported remarkable sensitivity values >95% for ovarian and liver
13
14 cancers. However, the reported sensitivity for breast cancer is rather low at 33%. The
15
16 novelty of our study is identifying circulating molecules that are associated with future
17
18 cancer risk and developing a pipeline to utilize these markers in generation of
19
20 biosensors based on our previous work to detect breast cancer risk. [41]
21
22
23
24
25

26 We used a combination of various statistical analysis methods to identify
27
28 biomarkers. Although, Student's t-test and/or random forest give some information
29
30 about the ability of a biomarker to discriminate between healthy and susceptible
31
32 patients, it alone is not sufficient. To identify the biomarkers with high classification
33
34 performance, we applied logistic regression. Area under curve (AUC) of receiver
35
36 operating characteristic (ROC) curves resulting from the classification operations on
37
38 these biomarkers are commonly used as an indicator for the discriminative capacity of a
39
40 single molecule or a set of molecules. Previously, logistic regression was performed on
41
42 predictors consisting of serum levels of several molecules, but authors did not report
43
44 any confidence interval for that AUC value and did not split the data into training and
45
46 test tests.[42] In another study, authors used Student' t-test and its non-parametric
47
48 equivalence (Mann-Whitney U-test) to find potential biomarkers, but the lower bound of
49
50 their reported confidence intervals was dramatically low, suggesting that those
51
52 biomarkers were not robust, and they also did not split the training and test sets.[43]
53
54
55
56
57
58
59
60
61
62
63
64
65

Several other studies have also used these methods to identify potential biomarkers but have not utilized a train-set split for their datasets.[17, 44-46] Training (Model building) and testing on the same dataset is not an ideal practice in machine learning as the model is likely to over-fit to the data. This approach results in artificially high predictive rates, in other words, *low generalizability*, which refers to poor applicability of the model to unseen data. The cross-validation that we employed in this study is a common approach to circumvent the problem of overfitting.

Conclusion

We identified biomarkers of breast cancer risk using metabolomics and protein profiling in plasma samples from healthy and susceptible individuals. Future studies are required to validate these markers in bigger data sets, to determine their role in breast tumorigenesis, develop liquid biopsy/biosensor-based approaches and move this information to clinic for early identification of breast cancer risk. In addition, further molecular studies in cell lines and animal models are required to show conclusively whether or not each or a combination of these markers can be utilized as indicators of breast cancer risk without having observable effects on breast cancer cells or can have other roles at the earlier stages of carcinogenesis. Overall, our analysis offers novel plasma biomarkers for further validation and functional characterization.

Abbreviations

Estrogen Receptor alpha, ER α

Progesterone Receptor, PgR

Human epidermal growth factor receptor, HER2

Breast cancer susceptibility 1 and 2 (BRCA1 and 2)

Circulating tumor cells, CTCs

Circulating tumor DNA, ctDNA

Institutional Review Board, IRB

University of Illinois, Urbana Champaign, UIUC

Gas chromatography-mass spectrometry, GC/MS

Area under curve, AUC

Receiver Operator Characteristic, ROC

Pearson's correlation coefficient, pCC

Logistic Regression, LR

Enzyme-linked immunosorbent assay, ELISA

Stem cell factor/KIT ligand, SCF/KITLG

Horse radish peroxidase, HRP

Deoxyribonucleic acid, DNA

Ribonucleic acid, RNA

Body mass index, BMI

Declarations

Ethical approval, informed consent to participate and publication: KTB studies were approved by the Indiana University Institutional Review Board (IRB protocol

number 1011003097 and 1607623663). All research was carried out in compliance with the Helsinki Declaration. Donors provided broad written consent for the use of their specimens in research. The consent document informed the donor that the donated specimens and medical data would be used for the general purpose of helping to determine how breast cancer develops and the exact laboratory experiments were unknown at the time of donation, and that proposals for use of the specimens would be reviewed and approved by a panel of independent researchers before specimens and/or data were released for research purposes.

Funding: This work was supported by grants from the University of Illinois, ACES Future interdisciplinary research explorations grant (ZME, data collection), Office of International Programs-Conrad Award (ZME, data collection), National Institute of Food and Agriculture, U.S. Department of Agriculture, award ILLU-698-909 (to ZME, data collection and analysis), UIUC Graduate College ASPIRE fellowship (ASC), Boğaziçi University research funds grant #12360 (to BA and HT, data analysis) and TÜBİTAK 2210-A National Scholarship Programme for MSc Students (KO).

Author Contributions: KO, BA, HT and ZME designed the study. KO performed the biomarker identification analysis. MP, NM and AMVS provided the samples. ASC prepared the samples for analysis and performed the experiments. KO, BA, HT, ASC and ZME wrote the manuscript. All authors have read and approved the manuscript.

Competing Interests: ZME has investigator-initiated grant from Karyopharm Therapeutics and is an coinventor on several patents entitled “Novel Compounds Which Activate Estrogen Receptors and Compositions and Methods of Using the Same. ZME was a PI on an investigator-initiated grant from Corteva Agrisciences and Pfizer Inc.

Data availability: All the data will be available from Komen Tissue Bank website upon acceptance of manuscript.

References

1. Siegel RL, Miller KD, Jemal A: **Cancer statistics, 2017**. *CA: A Cancer Journal for Clinicians* 2017, **67**(1):7-30.
2. Anderson BO, Yip C-H, Smith RA, Shyyan R, Sener SF, Eniu A, Carlson RW, Azavedo E, Harford J: **Guideline implementation for breast healthcare in low-income and middle-income countries**. *Cancer* 2008, **113**(8):2221-2243.
3. Coleman MP, Quaresma M, Berrino F, Lutz J-M, De Angelis R, Capocaccia R, Baili P, Rachet B, Gatta G, Hakulinen T *et al*: **Cancer survival in five continents: a worldwide population-based study (CONCORD)**. *The Lancet Oncology* 2008, **9**(8):730-756.
4. Li J, Shao Z: **Mammography screening in less developed countries**. *SpringerPlus* 2015, **4**:615.
5. da Costa Vieira RA, Biller G, Uemura G, Ruiz CA, Curado MP: **Breast cancer screening in developing countries**. *Clinics* 2017, **72**(4):244-253.
6. Dawson SJ, Duffy SW, Blows FM, Driver KE, Provenzano E, LeQuesne J, Greenberg DC, Pharoah P, Caldas C, Wishart GC: **Molecular characteristics of screen-detected vs symptomatic breast cancers and their impact on survival**. *Br J Cancer* 2009, **101**(8):1338-1344.
7. Iqbal J, Ginsburg O, Rochon PA, Sun P, Narod SA: **Differences in breast cancer stage at diagnosis and cancer-specific survival by race and ethnicity in the United States**. *JAMA* 2015, **313**(2):165-173.
8. Hammond MEH, Hayes DF, Dowsett M, Allred DC, Hagerty KL, Badve S, Fitzgibbons PL, Francis G, Goldstein NS, Hayes M *et al*: **American Society of Clinical Oncology/College of American Pathologists Guideline Recommendations for Immunohistochemical Testing of Estrogen and Progesterone Receptors in Breast Cancer**. *Journal of Clinical Oncology* 2010, **28**(16):2784-2795.
9. Andre F, Ismaila N, Stearns V: **Use of Biomarkers to Guide Decisions on Adjuvant Systemic Therapy for Women With Early-Stage Invasive Breast Cancer: ASCO Clinical Practice Guideline Update Summary**. *Journal of Oncology Practice* 2019, **15**(9):495-497.
10. Wolff AC, Hammond MEH, Allison KH, Harvey BE, Mangu PB, Bartlett JMS, Bilous M, Ellis IO, Fitzgibbons P, Hanna W *et al*: **Human Epidermal Growth Factor Receptor 2 Testing in Breast Cancer: American Society of Clinical Oncology/College of American Pathologists Clinical Practice Guideline Focused Update**. *Journal of Clinical Oncology* 2018, **36**(20):2105-2122.
11. Van Poznak C, Somerfield MR, Bast RC, Cristofanilli M, Goetz MP, Gonzalez-Angulo AM, Hicks DG, Hill EG, Liu MC, Lucas W *et al*: **Use of Biomarkers to Guide Decisions on Systemic Therapy for Women With Metastatic Breast Cancer: American Society of Clinical Oncology Clinical Practice Guideline**. *Journal of Clinical Oncology* 2015, **33**(24):2695-2704.
12. Goossens N, Nakagawa S, Sun X, Hoshida Y: **Cancer biomarker discovery and validation**. *Translational cancer research* 2015, **4**(3):256-269.
13. O'Leary B, Hrebien S, Morden JP, Beaney M, Fribbens C, Huang X, Liu Y, Bartlett CH, Koehler M, Cristofanilli M *et al*: **Early circulating tumor DNA dynamics and clonal selection with palbociclib and fulvestrant for breast cancer**. *Nature Communications* 2018, **9**(1):896.
14. Coombes RC, Page K, Salari R, Hastings RK, Armstrong AC, Ahmed S, Ali S, Cleator SJ, Kenny LM, Stebbing J *et al*: **Personalized detection of circulating tumor DNA antedates breast cancer metastatic recurrence**. *Clinical Cancer Research* 2019:clincanres.3663.2018.

15. Merker JD, Oxnard GR, Compton C, Diehn M, Hurley P, Lazar AJ, Lindeman N, Lockwood CM, Rai AJ, Schilsky RL *et al*: **Circulating Tumor DNA Analysis in Patients With Cancer: American Society of Clinical Oncology and College of American Pathologists Joint Review**. *Journal of Clinical Oncology* 2018, **36**(16):1631-1641.
16. Madak-Erdogan Z, Band S, Zhao YC, Smith BP, Kulkoyluoglu-Cotul E, Zuo Q, Santaliz Casiano A, Wrobel K, Rossi G, Smith RL *et al*: **Free fatty acids rewire cancer metabolism in obesity-associated breast cancer via estrogen receptor and mTOR signaling**. *Cancer Res* 2019.
17. Kazarian A, Blyuss O, Metodieva G, Gentry-Maharaj A, Ryan A, Kiseleva EM, Prytomanova OM, Jacobs IJ, Widschwendter M, Menon U *et al*: **Testing breast cancer serum biomarkers for early detection and prognosis in pre-diagnosis samples**. *British Journal of Cancer* 2017, **116**(4):501-508.
18. Hines S, Organ C, Kornstein M, Krystal G: **Coexpression of the c-kit and stem cell factor genes in breast carcinomas**. *Cell Growth Differ* 1995, **6**(6):769-779.
19. Harrell JC, Shroka TM, Jacobsen BM: **Estrogen induces c-Kit and an aggressive phenotype in a model of invasive lobular breast cancer**. *Oncogenesis* 2017, **6**(11):396.
20. Opyrchal M, Gil M, Salisbury JL, Goetz MP, Suman V, Degnim A, McCubrey J, Haddad T, Iankov I, Kurokawa CB *et al*: **Molecular targeting of the Aurora-A/SMAD5 oncogenic axis restores chemosensitivity in human breast cancer cells**. *Oncotarget* 2017, **8**(53):91803-91816.
21. Opyrchal M, Iankov I, Ingle JN, Salisbury JL, Galanis E, D'Assoro A: **SMAD5 expression and inhibition of the mitotic kinase aurora-A on sensitivity of breast cancer cells to chemotherapy**. *Journal of Clinical Oncology* 2013, **31**(15_suppl):e13516-e13516.
22. Ghassemi S, Vejdovszky K, Sahin E, Ratzinger L, Schelch K, Mohr T, Peter-Vörösmarty B, Brankovic J, Lackner A, Leopoldi A *et al*: **FGF5 is expressed in melanoma and enhances malignancy in vitro and in vivo**. *Oncotarget* 2017, **8**(50):87750-87762.
23. Huang Y, Wang H, Yang Y: **Expression of Fibroblast Growth Factor 5 (FGF5) and Its Influence on Survival of Breast Cancer Patients**. *Med Sci Monit* 2018, **24**:3524-3530.
24. Ates O, Gedik E, Sunar V, Altundag K: **Serum endocan level and its prognostic significance in breast cancer patients**. *Journal of Oncological Sciences* 2018, **4**(1):15-18.
25. Bębenek M, Duś D, Koźlak J: **Prognostic value of the Fas/Fas ligand system in breast cancer**. *Contemporary Oncology* 2013, **17**(2):120-122.
26. Mor G, Kohen F, Garcia-Velasco J, Nilsen J, Brown W, Song J, Naftolin F: **Regulation of Fas ligand expression in breast cancer cells by estrogen: functional differences between estradiol and tamoxifen**. *The Journal of Steroid Biochemistry and Molecular Biology* 2000, **73**(5):185-194.
27. Toss M, Miligy I, Gorringer K, Mittal K, Aneja R, Ellis I, Green A, Rakha E: **Prognostic significance of cathepsin V (CTSV/CTSL2) in breast ductal carcinoma in situ**. *Journal of Clinical Pathology* 2019:jclinpath-2019-205939.
28. Benson CS, Babu SD, Radhakrishna S, Selvamurugan N, Sankar BR: **Expression of Matrix Metalloproteinases in Human Breast Cancer Tissues**. *Disease markers* 2013, **34**(6):395-405.
29. Zhang G, Miyake M, Lawton A, Goodison S, Rosser CJ: **Matrix metalloproteinase-10 promotes tumor progression through regulation of angiogenic and apoptotic pathways in cervical tumors**. *BMC Cancer* 2014, **14**(1):310.
30. Cheng T, Wei R, Jiang G, Zhou Y, Lv M, Dai Y, Yuan Y, Luo D, Ma D, Li F *et al*: **XPNPEP2 is overexpressed in cervical cancer and promotes cervical cancer metastasis**. *Tumor Biology* 2017, **39**(7):1010428317717122.
31. Knott ME, Ranuncolo SM, Nuñez M, Armanasco E, Puricelli LI, De Lorenzo MS: **Abstract 1577: Levels of Fibroblast Growth Factor 21 (FGF21) in serum as diagnostic biomarker in patients with breast cancer**. *Cancer Research* 2015, **75**(15 Supplement):1577-1577.

32. Aukes K, Forsman C, Brady NJ, Astleford K, Blixt N, Sachdev D, Jensen ED, Mansky KC, Schwertfeger KL: **Breast cancer cell-derived fibroblast growth factors enhance osteoclast activity and contribute to the formation of metastatic lesions.** *PLOS ONE* 2017, **12**(10):e0185736.
33. Matsui A, Yokoo H, Negishi Y, Endo-Takahashi Y, Chun NAL, Kadouchi I, Suzuki R, Maruyama K, Aramaki Y, Semba K *et al*: **CXCL17 expression by tumor cells recruits CD11b+Gr1 high F4/80-cells and promotes tumor progression.** *PloS one* 2012, **7**(8):e44080-e44080.
34. Ben-Baruch A, Xu L, Young PR, Bengali K, Oppenheim JJ, Wang JM: **Monocyte Chemotactic Protein-3 (MCP3) Interacts with Multiple Leukocyte Receptors: C-C CKR1, A RECEPTOR FOR MACROPHAGE INFLAMMATORY PROTEIN-1 α /RANTES, IS ALSO A FUNCTIONAL RECEPTOR FOR MCP3.** *Journal of Biological Chemistry* 1995, **270**(38):22123-22128.
35. Sun H, Xu J, Huang Q, Huang M, Li K, Qu K, Wen H, Lin R, Zheng M, Wei H *et al*: **Reduced CD160 Expression Contributes to Impaired NK-cell Function and Poor Clinical Outcomes in Patients with HCC.** *Cancer Research* 2018, **78**(23):6581-6593.
36. Huang Y, Yu X, Wang L, Zhou S, Sun J, Feng N, Nie S, Wu J, Gao F, Fei B *et al*: **Four Genetic Polymorphisms of Lymphotoxin-Alpha Gene and Cancer Risk: A Systematic Review and Meta-Analysis.** *PLOS ONE* 2013, **8**(12):e82519.
37. Kohaar I, Tiwari P, Kumar R, Nasare V, Thakur N, C Das B, Bharadwaj M: **Association of single nucleotide polymorphisms (SNPs) in TNF-LTA locus with breast cancer risk in Indian population,** vol. 114; 2008.
38. Alix-Panabières C, Pantel K: **Circulating Tumor Cells: Liquid Biopsy of Cancer.** *Clinical Chemistry* 2013, **59**(1):110-118.
39. Heitzer E, Ulz P, Geigl JB: **Circulating Tumor DNA as a Liquid Biopsy for Cancer.** *Clinical Chemistry* 2015, **61**(1):112-123.
40. Cohen JD, Li L, Wang Y, Thoburn C, Afsari B, Danilova L, Douville C, Javed AA, Wong F, Mattox A *et al*: **Detection and localization of surgically resectable cancers with a multi-analyte blood test.** *Science (New York, NY)* 2018, **359**(6378):926-930.
41. Sarangapani K, Torun H, Finkler O, Zhu C, Degertekin L: **Membrane-based actuation for high-speed single molecule force spectroscopy studies using AFM.** *European Biophysics Journal* 2010, **39**(8):1219-1227.
42. Hwa H-L, Kuo W-H, Chang L-Y, Wang M-Y, Tung T-H, Chang K-J, Hsieh F-J: **Prediction of breast cancer and lymph node metastatic status with tumour markers using logistic regression models.** *Journal of Evaluation in Clinical Practice* 2008, **14**(2):275-280.
43. Santillán-Benítez JG, Mendieta-Zerón H, Gómez-Oliván LM, Torres-Juárez JJ, González-Bañales JM, Hernández-Peña LV, Ordóñez-Quiroz A: **The Tetrad BMI, Leptin, Leptin/Adiponectin (L/A) Ratio and CA 15-3 are Reliable Biomarkers of Breast Cancer.** *Journal of Clinical Laboratory Analysis* 2013, **27**(1):12-20.
44. Dalamaga M, Sotiropoulos G, Karmaniolas K, Pelekanos N, Papadavid E, Lekka A: **Serum resistin: A biomarker of breast cancer in postmenopausal women? Association with clinicopathological characteristics, tumor markers, inflammatory and metabolic parameters.** *Clinical Biochemistry* 2013, **46**(7):584-590.
45. Assiri AMA, Kamel HFM, Hassanien MFR: **Resistin, Visfatin, Adiponectin, and Leptin: Risk of Breast Cancer in Pre- and Postmenopausal Saudi Females and Their Possible Diagnostic and Predictive Implications as Novel Biomarkers.** *Disease Markers* 2015, **2015**:253519.
46. Provatopoulou X, Georgiou GP, Kalogera E, Kalles V, Matiatou MA, Papapanagiotou I, Sagkriotis A, Zografos GC, Gounaris A: **Serum irisin levels are lower in patients with breast cancer: association with disease diagnosis and tumor characteristics.** *BMC Cancer* 2015, **15**:898.

Table and Figure legends

Table 1. The p-values of the paired t-test analysis for each pair of molecules among the top 20 molecules ranked by applying the Student's t-test to all 181 molecules. The paired t-test assesses the pairwise dependencies of the most discriminative 20 molecules. Pairs with $p > 0.05$ show strong dependency within that pair.

Table 2. The p-values of the paired t-test analysis for each pair of molecules among the top 20 molecules ranked by applying the random forest to all 181 molecules. The paired t-test assesses the pairwise dependencies of the most discriminative 20 molecules. Pairs with $p > 0.05$ show strong dependency within that pair.

Table 3. Pearson's correlation coefficient between each pair of 20 most important molecules ranked by their Student's t-test p-values. Pairs with $pCC > 0.5$ show strong dependency within that pair.

Table 4. Pearson's correlation coefficient between each pair of 20 most important molecules ranked by their importance calculated by random forest. Pairs with $pCC > 0.5$ show strong dependency within that pair.

Table 5. Ranking of molecules identified by initial Student's t-test for each consequent feature elimination method. The "Student's t-test" column lists the top ranking most discriminative 20 molecules among 181 molecules. The "Manual Selection by LR", "Paired t-test" and "Correlation Analysis" columns list the molecules

selected from these 20 top molecules by applying the iterative molecule elimination procedures using manual selection by LR based on classification performance, paired t-test and correlation analysis, respectively, as described in Statistical Analysis.

Table 6. Ranking of molecules identified by initial random forest method for each consequent feature elimination method. The “Random Forest” column lists the top ranking most discriminative 20 molecules among 181 molecules. The “Manual Selection by LR”, “Paired t-test” and “Correlation Analysis” columns list the molecules selected from these 20 top molecules by applying the iterative molecule elimination procedures using manual selection by LR based on classification performance, paired t-test and correlation analysis, respectively, as described in Statistical Analysis.

Figure 1. Identification and performance assessment of circulating factor signatures for future breast cancer risk assessment using Student’s t-test as initial feature selection method (A) Levels of top 20 molecules identified by Student’s t-test in 47 healthy (red) and 49 susceptible (green) individuals using OLINK analysis. Z-scores were not log transformed or centered. Unsupervised hierarchical clustering was performed using Cluster 3 software for Z-scores of molecule concentrations with uncentered correlation as similarity metric and average linkage as clustering method. Data are visualized using Java Tree view software. In the lower panel, each column represents an individual and each row represents a molecule, with elevated levels in red, reduced levels in blue, and mean control levels in white. Bar indicates the coloring for Z-scores of molecule concentrations. (B) LR classification performances (AUC values) using the top-ranking N (1-20) molecules, ranked by their p-values in Table 5, and the ROC curves of every AUC value where the bold black line indicates ROC curve

of the best performing (the highest AUC value) molecule set. (C) LR classification performances (AUC values) using the top-ranking N (1-13) molecules selected manually by considering the LR classification performance given in Figure 1B and the ROC curves of every AUC value where the bold black line indicates ROC curve of the best performing (the highest AUC value) molecule set. (D) LR classification performances (AUC values) using the top-ranking molecules selected from the list of 20 molecules, ranked by Student's t-test in Table 5, by iterative elimination using pair-wise Pearson correlation coefficients of features in Table 3 ($|pCC|=0.5$ is the significance threshold). The ROC curves of every AUC value where the bold black line indicates ROC curve of the best performing (the highest AUC value) molecule set. (E) LR classification performances (AUC values) using the top-ranking N (1-2) molecules selected from the list of 20 molecules, ranked by Student's t-test in Table 5, by iterative elimination using paired t-test p-values of features in Table 1 ($p=0.05$ is the significance threshold). The ROC curves of every AUC value where the bold black line indicates ROC curve of the best performing (the highest AUC value) molecule set.

Figure 2. Identification and performance assessment of circulating factor signatures for future breast cancer risk assessment using random forest as initial feature selection method (A) Levels of top 20 molecules identified by random forest method in 47 healthy (red) and 49 susceptible (green) individuals using OLINK analysis. Z-Scores were not log transformed or centered. Unsupervised hierarchical clustering was performed using Cluster 3 software for Z-scores of molecule concentrations with uncentered correlation as similarity metric and average linkage as clustering method. Data are visualized using Java Tree view software. In the lower panel, each column

represents an individual and each row represents a molecule, with elevated levels in red, reduced levels in blue, and mean control levels in white. Bar indicates the coloring for Z-scores of molecule concentrations. (B) LR classification performances (AUC values) using the top-ranking N (1-20) molecules, ranked by their feature importance values (computed by random forest) in Table 6, and the ROC curves of every AUC value where the bold black line indicates ROC curve of the best performing (the highest AUC value) molecule set. (C) LR classification performances (AUC values) using the top-ranking molecules selected manually by considering the LR classification performance given in Figure 2B and the ROC curves of every AUC value where the bold black line indicates ROC curve of the best performing (the highest AUC value) molecule set. (D) LR classification performances (AUC values) using the top-ranking N (1-11) molecules selected from the list of 20 molecules, ranked by random forest in Table 6, by iterative elimination using pair-wise Pearson correlation coefficients of features in Table 4 ($|pCC|=0.5$ is the significance threshold). The ROC curves of every AUC value where the bold black line indicates ROC curve of the best performing (the highest AUC value) molecule set. (E) LR classification performances (AUC values) using the top-ranking N (1-2) molecules selected from the list of 20 molecules, ranked by random forest in Table 6, by iterative elimination using paired t-test p-values of features in Table 2 ($p=0.05$ is the significance threshold). The ROC curves of every AUC value where the bold black line indicates ROC curve of the best performing (the highest AUC value) molecule set.

Figure 3. A. Validation of biomarker identification. (A) Levels of 13 molecules identified by Student's t-test followed by manual selection in 47 healthy (red) and 49

susceptible (green) individuals using OLINK analysis. Z-Scores were not log transformed or centered. Unsupervised hierarchical clustering was performed using Cluster 3 software for Z-scores of molecule concentrations with uncentered correlation as similarity metric and average linkage as clustering method. Data are visualized using Java Tree view software. In the lower panel, each column represents an individual and each row represents a molecule, with elevated levels in red, reduced levels in blue, and mean control levels in white. Bar indicates the coloring for Z-scores of molecule concentrations. (B) Changes in the levels of 12 of 13 signature molecules in 47 healthy and 49 susceptible individuals. Anderson-Darling and Kolmogorov-Smirnov tests for normality was used. If the dataset didn't pass the normality test, non-parametric Mann-Whitney test was used to assess if level of a molecule is statistically significantly different in plasma from healthy vs. susceptible individuals (molecules with *). Otherwise, unpaired t-test was used to assess if level of a molecule is statistically significantly different in plasma from healthy vs. susceptible individuals. All data points are plotted. P-values are indicted on the graphs. (C) Level of SCF/KITLG in 47 healthy and 49 susceptible individuals. Anderson-Darling and Kolmogorov-Smirnov tests for normality was used. Non-parametric Mann-Whitney test was used to assess if level of a molecule is statistically significantly different in plasma from healthy vs. susceptible individuals. All data points are plotted (as histogram on the left side and as box-whiskers graph on right side). P-values are indicted on the graph. (D) Results from 3C are independently validated using ELISA assay using all the samples. Level of identified biomarkers in human plasma samples were compared using unpaired t-test. P-value is reported on the graph. Values from all the samples are presented.

Table 1

Table 1	SCF	MAD HOMO	FGF-5	FASLG	MMP-10	PPY
SCF	0	0.6963	0.0004	0.0006	0.3286	0.0017
MAD HOMOLOG 5	0.6963	0	0.0013	0.0014	0.6643	0.0038
FGF-5	0.0004	0.0013	0	0.6618	0.002	0.7183
FASLG	0.0006	0.0014	0.6618	0	0.0018	0.9554
MMP-10	0.3286	0.6643	0.002	0.0018	0	0.0033
PPY	0.0017	0.0038	0.7183	0.9554	0.0033	0
XPNPEP2	0.361	0.5633	0.0018	0.0015	0.9115	0.0054
FGF-21	0.003	0.0026	0.6703	0.9738	0.0006	0.9356
CXL17	0.2659	0.4489	0.0017	0.0034	0.8086	0.0018
MCP-3	0.189	0.2985	0.0027	0.002	0.542	0.0055
ESM-1	0.0003	0.0006	0.6315	0.9137	0.0073	0.8852
HK11	0.1552	0.3257	0.0064	0.0015	0.5825	0.0063
TRAIL	0.1309	0.5298	0.0033	0.0039	0.8731	0.0096
FGF-BP1	0.1303	0.2913	0.01	0.0112	0.5308	0.0166
EN-RAGE	0.1191	0.354	0.0027	0.005	0.6119	0.0124
C15:0	0.1324	0.2454	0.0091	0.0067	0.417	0.0109
TNFB	0.0029	0.008	0.4087	0.6151	0.0071	0.6315
CTSV	0.0045	0.0062	0.4174	0.548	0.0154	0.6427
ADA	0.001	0.0083	0.6858	0.9595	0.0169	0.93
CD160	0.0015	0.007	0.4709	0.6497	0.0086	0.6801

XPNPEP2	FGF-21	CXL17	MCP-3	ESM-1	HK11	TRAIL	FGF-BP1
0.361	0.003	0.2659	0.189	0.0003	0.1552	0.1309	0.1303
0.5633	0.0026	0.4489	0.2985	0.0006	0.3257	0.5298	0.2913
0.0018	0.6703	0.0017	0.0027	0.6315	0.0064	0.0033	0.01
0.0015	0.9738	0.0034	0.002	0.9137	0.0015	0.0039	0.0112
0.9115	0.0006	0.8086	0.542	0.0073	0.5825	0.8731	0.5308
0.0054	0.9356	0.0018	0.0055	0.8852	0.0063	0.0096	0.0166
0	0.0034	0.9172	0.66	0.0125	0.6878	0.9708	0.6209
0.0034	0	0.0014	0.0029	0.9489	0.006	0.0063	0.013
0.9172	0.0014	0	0.7384	0.0038	0.7438	0.9536	0.6714
0.66	0.0029	0.7384	0	0.0089	0.9712	0.7017	0.9032
0.0125	0.9489	0.0038	0.0089	0	0.0032	0.0061	0.013
0.6878	0.006	0.7438	0.9712	0.0032	0	0.7392	0.8734
0.9708	0.0063	0.9536	0.7017	0.0061	0.7392	0	0.6469
0.6209	0.013	0.6714	0.9032	0.013	0.8734	0.6469	0
0.6971	0.0134	0.7648	0.9911	0.0112	0.9815	0.6816	0.9125
0.4563	0.0096	0.5413	0.788	0.0289	0.7421	0.5747	0.8767
0.0108	0.7003	0.0172	0.0183	0.7453	0.0202	0.0092	0.0416
0.0055	0.6644	0.026	0.017	0.7257	0.0234	0.0274	0.0442
0.0188	0.9805	0.019	0.0211	0.9715	0.037	0.0071	0.0429
0.0214	0.7499	0.0126	0.0241	0.7768	0.0087	0.0078	0.0441

EN-RAGE	C15:0	TNFB	CTSV	ADA	CD160
0.1191	0.1324	0.0029	0.0045	0.001	0.0015
0.354	0.2454	0.008	0.0062	0.0083	0.007
0.0027	0.0091	0.4087	0.4174	0.6858	0.4709
0.005	0.0067	0.6151	0.548	0.9595	0.6497
0.6119	0.417	0.0071	0.0154	0.0169	0.0086
0.0124	0.0109	0.6315	0.6427	0.93	0.6801
0.6971	0.4563	0.0108	0.0055	0.0188	0.0214
0.0134	0.0096	0.7003	0.6644	0.9805	0.7499
0.7648	0.5413	0.0172	0.026	0.019	0.0126
0.9911	0.788	0.0183	0.017	0.0211	0.0241
0.0112	0.0289	0.7453	0.7257	0.9715	0.7768
0.9815	0.7421	0.0202	0.0234	0.037	0.0087
0.6816	0.5747	0.0092	0.0274	0.0071	0.0078
0.9125	0.8767	0.0416	0.0442	0.0429	0.0441
0	0.8043	0.0374	0.025	0.0115	0.0327
0.8043	0	0.0409	0.0464	0.0725	0.0502
0.0374	0.0409	0	0.969	0.7488	0.9447
0.025	0.0464	0.969	0	0.7132	0.9266
0.0115	0.0725	0.7488	0.7132	0	0.8146
0.0327	0.0502	0.9447	0.9266	0.8146	0

Table 2

Table 2	SCF	MAD	HOMOL	PPV	FASLG	FGF-5	CXCL1	MMP-10	XPNPEP2	ESM-1	PHOSPHORIC
SCF	0	0.6963		0.0017	0.0006	0.0004				0.0003	0.0859
MAD	HOMOL	0.6963	0	0.0038	0.0014	0.0013	0.0153	0.6643	0.5633	0.0006	0.1622
PPV	0.0017	0.0038		0	0.9554	0.7183	0.4656	0.0033	0.0054	0.8852	0.0413
FASLG	0.0006	0.0014		0.9554	0	0.6618	0.4225	0.0018	0.0015	0.9137	0.0376
FGF-5	0.0004	0.0013		0.7183	0		0.3008	0.002	0.0018	0.6315	0.0209
CXCL1	0.0023	0.0153		0.4656	0.4225	0		0.0115	0.0376	0.5516	0.2144
MMP-10	0.3286	0.6643		0.0033	0.0018	0.002	0.0115	0	0.9115	0.0073	0.3345
XPNPEP2	0.361	0.5633		0.0054	0.0015	0.0018	0.0376	0.9115	0	0.0125	0.3469
ESM-1	0.0003	0.0006		0.8852	0.9137	0.6315	0.5516	0.0073	0.0125	0	0.0593
PHOSPHORIC	0.0859	0.1622		0.0413	0.0376	0.0209	0.2144	0.3345	0.3469	0.0593	0
PD-L1	0.0333	0.146		0.0721	0.0732	0.0378	0.2567	0.1225	0.2814	0.114	0.829
EPHA2	0.0073	0.0305		0.1935	0.142	0.1085	0.5558	0.0486	0.1044	0.1401	0.4717
FLT3L	0.0056	0.0442		0.2713	0.2782	0.1423	0.6492	0.0419	0.0921	0.3336	0.446
4E-BP1	0.1307	0.2792		0.0397	0.0277	0.0145	0.1072	0.5293	0.5733	0.0335	0.7561
TRAIL	0.1309	0.5298		0.0096	0.0039	0.0033	0.0076	0.8731	0.9708	0.0061	0.422
MCP-1	0.0148	0.0506		0.1833	0.166	0.0857	0.4841	0.0542	0.1341	0.2512	0.517
TLR3	0.0432	0.0586		0.2008	0.1646	0.0983	0.5033	0.1309	0.1636	0.2426	0.5868
CD27	0.0116	0.0252		0.2943	0.2553	0.1787	0.691	0.0226	0.0719	0.3007	0.3887
FGF-8P1	0.1303	0.2913		0.0166	0.0112	0.01	0.0583	0.5308	0.6209	0.013	0.6436
HK14	0.1215	0.2387		0.0205	0.0051	0.0147	0.0429	0.4234	0.519	0.0129	0.762

PD-L1	EPHA2	FLT3L	4E-BP1	TRAIL	MCP-1	TLR3	CD27	FGF-BP1	HK14
0.0333	0.0073	0.0056	0.1307	0.1309	0.0148	0.0432	0.0116	0.1303	0.1215
0.146	0.0305	0.0442	0.2792	0.5298	0.0506	0.0586	0.0252	0.2913	0.2387
0.0721	0.1935	0.2713	0.0397	0.0096	0.1833	0.2008	0.2943	0.0166	0.0205
0.0732	0.142	0.2782	0.0277	0.0039	0.166	0.1646	0.2553	0.0112	0.0051
0.0378	0.1085	0.1423	0.0145	0.0033	0.0857	0.0983	0.1787	0.01	0.0147
0.2567	0.5558	0.6492	0.1072	0.0076	0.4841	0.5033	0.691	0.0583	0.0429
0.1225	0.0486	0.0419	0.5293	0.8731	0.0542	0.1309	0.0226	0.5308	0.4234
0.2814	0.1044	0.0921	0.5733	0.9708	0.1341	0.1636	0.0719	0.6209	0.519
0.114	0.1401	0.3336	0.0335	0.0061	0.2512	0.2426	0.3007	0.013	0.0129
0.829	0.4717	0.446	0.7561	0.422	0.517	0.5868	0.3887	0.6436	0.762
0	0.5517	0.4038	0.5749	0.1934	0.5838	0.7439	0.415	0.5413	0.547
0.5517	0	0.8913	0.3318	0.0648	0.9749	0.9147	0.7529	0.2438	0.2666
0.4038	0.8913	0	0.2581	0.0333	0.8507	0.8279	0.94	0.2181	0.2617
0.5749	0.3318	0.2581	0	0.5541	0.1918	0.4236	0.2596	0.924	0.9861
0.1934	0.0648	0.0333	0.5541	0	0.066	0.1962	0.0647	0.6469	0.5674
0.5838	0.9749	0.8507	0.1918	0.066	0	0.9421	0.8181	0.2741	0.3109
0.7439	0.9147	0.8279	0.4236	0.1962	0.9421	0	0.7566	0.3418	0.3587
0.415	0.7529	0.94	0.2596	0.0647	0.8181	0.7566	0	0.1532	0.1312
0.5413	0.2438	0.2181	0.924	0.6469	0.2741	0.3418	0.1532	0	0.8822
0.547	0.2666	0.2617	0.9861	0.5674	0.3109	0.3587	0.1312	0.8822	0

Table 3

Table 3	SCF	MAD HOMO	FGF-5	FASLG	MMP-10	PPY	XPNPEP2
SCF	1	0.136	0.129	0.076	0.318	-0.088	0.03
MAD HOMO	0.136	1	0.057	0.064	-0.048	-0.119	0.148
FGF-5	0.129	0.057	1	0.21	0.047	0.15	0.123
FASLG	0.076	0.064	0.21	1	0.079	0.169	0.179
MMP-10	0.318	-0.048	0.047	0.079	1	-0.002	0.156
PPY	-0.088	-0.119	0.15	0.169	-0.002	1	-0.047
XPNPEP2	0.03	0.148	0.123	0.179	0.156	-0.047	1
FGF-21	-0.189	0.021	0.104	0.192	0.32	0.085	0.11
CXL17	0.215	0.323	0.151	0.037	0.323	0.199	0.168
MCP-3	0.03	0.217	0.116	0.165	0.281	0.012	0.15
ESM-1	0.265	0.305	0.067	0.149	-0.087	-0.006	-0.169
HK11	0.261	0.208	-0.048	0.246	0.264	0.022	0.188
TRAIL	0.66	0.244	0.113	0.127	0.294	-0.065	-0.03
FGF-BP1	0.228	0.118	-0.073	-0.068	0.076	-0.147	0.004
EN-RAGE	0.415	0.132	0.233	0.154	0.143	-0.034	0.118
C15:0	0	0.017	-0.027	0.075	0.097	0.003	0.227
TNFB	0.073	0.012	0.184	0.285	0.159	0.089	0.138
CTSV	-0.024	0.083	0.049	0.384	-0.042	-0.175	0.281
ADA	0.312	0.089	0.219	0.12	0.043	0.035	0.083
CD160	0.223	0.082	0.161	0.446	0.158	0.173	-0.018

FGF-21	CXL17	MCP-3	ESM-1	HK11	TRAIL	FGF-BP1	EN-RAGE
-0.189	0.215	0.03	0.265	0.261	0.66	0.228	0.415
0.021	0.323	0.217	0.305	0.208	0.244	0.118	0.132
0.104	0.151	0.116	0.067	-0.048	0.113	-0.073	0.233
0.192	0.037	0.165	0.149	0.246	0.127	-0.068	0.154
0.32	0.323	0.281	-0.087	0.264	0.294	0.076	0.143
0.085	0.199	0.012	-0.006	0.022	-0.065	-0.147	-0.034
0.11	0.168	0.15	-0.169	0.188	-0.03	0.004	0.118
1	0.276	0.2	-0.124	0.092	0.101	0.003	0.005
0.276	1	0.129	0.164	0.32	0.198	0.092	0.141
0.2	0.129	1	0.02	0.068	0.214	0.112	0.319
-0.124	0.164	0.02	1	0.269	0.16	0.072	0.117
0.092	0.32	0.068	0.269	1	0.176	0.187	0.047
0.101	0.198	0.214	0.16	0.176	1	0.134	0.435
0.003	0.092	0.112	0.072	0.187	0.134	1	-0.055
0.005	0.141	0.319	0.117	0.047	0.435	-0.055	1
0.103	0.131	-0.112	-0.187	0.13	-0.066	0.008	-0.136
0.016	0.042	0.11	0.075	0.112	0.28	-0.017	0.019
0.038	-0.096	0.133	-0.026	0.071	-0.009	-0.034	0.165
0.065	0.119	0.202	0.233	0.007	0.379	0.062	0.41
0.097	0.17	0.079	0.299	0.352	0.333	0.006	0.107

C15:0	TNFB	CTSV	ADA	CD160
0	0.073	-0.024	0.312	0.223
0.017	0.012	0.083	0.089	0.082
-0.027	0.184	0.049	0.219	0.161
0.075	0.285	0.384	0.12	0.446
0.097	0.159	-0.042	0.043	0.158
0.003	0.089	-0.175	0.035	0.173
0.227	0.138	0.281	0.083	-0.018
0.103	0.016	0.038	0.065	0.097
0.131	0.042	-0.096	0.119	0.17
-0.112	0.11	0.133	0.202	0.079
-0.187	0.075	-0.026	0.233	0.299
0.13	0.112	0.071	0.007	0.352
-0.066	0.28	-0.009	0.379	0.333
0.008	-0.017	-0.034	0.062	0.006
-0.136	0.019	0.165	0.41	0.107
1	0.035	-0.012	-0.162	-0.001
0.035	1	0.048	0.204	0.275
-0.012	0.048	1	0.25	-0.087
-0.162	0.204	0.25	1	0.054
-0.001	0.275	-0.087	0.054	1

Table 4

Table 4	SCF	MAD HOMO	PPY	FASLG	FGF-5	CXCL1	MMP-10
SCF	1	0.136	-0.088	0.076	0.129	0.227	0.318
MAD HOMO	0.136	1	-0.119	0.064	0.057	-0.029	-0.048
PPY	-0.088	-0.119	1	0.169	0.15	-0.038	-0.002
FASLG	0.076	0.064	0.169	1	0.21	0.19	0.079
FGF-5	0.129	0.057	0.15	0.21	1	0.044	0.047
CXCL1	0.227	-0.029	-0.038	0.19	0.044	1	0.195
MMP-10	0.318	-0.048	-0.002	0.079	0.047	0.195	1
XPNPEP2	0.03	0.148	-0.047	0.179	0.123	-0.094	0.156
ESM-1	0.265	0.305	-0.006	0.149	0.067	0.01	-0.087
PHOSPHORIC	0.05	0.067	-0.016	-0.011	0.09	-0.158	-0.111
PD-L1	0.363	0.061	0.235	0.229	0.28	0.317	0.499
EPHA2	0.382	0.271	0.148	0.347	0.218	0.197	0.352
FLT3L	0.395	0.105	0.102	0.098	0.253	0.263	0.348
4E-BP1	0.297	0.231	-0.014	0.107	0.18	0.202	0.101
TRAIL	0.66	0.244	-0.065	0.127	0.113	0.439	0.294
MCP-1	0.277	0.152	0.184	0.262	0.319	0.402	0.369
TLR3	-0.011	0.15	-0.017	0.149	0.178	0.126	0.047
CD27	0.131	0.133	-0.011	0.159	0.081	0.235	0.343
FGF-BP1	0.228	0.118	-0.147	-0.068	-0.073	0.066	0.076
HK14	0.159	0.157	-0.072	0.246	-0.068	0.319	0.209

XPNPEP2	ESM-1	PHOSPHORIC PD-L1	EPHA2	FLT3L	4E-BP1	TRAIL	
0.03	0.265	0.05	0.363	0.382	0.395	0.297	0.66
0.148	0.305	0.067	0.061	0.271	0.105	0.231	0.244
-0.047	-0.006	-0.016	0.235	0.148	0.102	-0.014	-0.065
0.179	0.149	-0.011	0.229	0.347	0.098	0.107	0.127
0.123	0.067	0.09	0.28	0.218	0.253	0.18	0.113
-0.094	0.01	-0.158	0.317	0.197	0.263	0.202	0.439
0.156	-0.087	-0.111	0.499	0.352	0.348	0.101	0.294
1	-0.169	0.063	0.119	0.151	0.145	0.2	-0.03
-0.169	1	-0.015	0.166	0.489	0.117	0.179	0.16
0.063	-0.015	1	-0.046	0.073	-0.066	0.082	-0.058
0.119	0.166	-0.046	1	0.525	0.635	0.316	0.484
0.151	0.489	0.073	0.525	1	0.358	0.165	0.431
0.145	0.117	-0.066	0.635	0.358	1	0.27	0.534
0.2	0.179	0.082	0.316	0.165	0.27	1	0.422
-0.03	0.16	-0.058	0.484	0.431	0.534	0.422	1
0.071	0.136	-0.015	0.519	0.368	0.507	0.579	0.466
0.02	0.029	-0.138	0.083	0.17	0.128	0.036	0.01
0.048	0.23	-0.166	0.452	0.642	0.159	0.054	0.217
0.004	0.072	-0.009	-0.053	0.064	0.043	0.067	0.134
0.121	0.196	-0.196	0.296	0.195	0.075	0.12	0.13

MCP-1	TLR3	CD27	FGF-BP1	HK14
0.277	-0.011	0.131	0.228	0.159
0.152	0.15	0.133	0.118	0.157
0.184	-0.017	-0.011	-0.147	-0.072
0.262	0.149	0.159	-0.068	0.246
0.319	0.178	0.081	-0.073	-0.068
0.402	0.126	0.235	0.066	0.319
0.369	0.047	0.343	0.076	0.209
0.071	0.02	0.048	0.004	0.121
0.136	0.029	0.23	0.072	0.196
-0.015	-0.138	-0.166	-0.009	-0.196
0.519	0.083	0.452	-0.053	0.296
0.368	0.17	0.642	0.064	0.195
0.507	0.128	0.159	0.043	0.075
0.579	0.036	0.054	0.067	0.12
0.466	0.01	0.217	0.134	0.13
1	0.042	0.155	0.03	0.121
0.042	1	0.138	-0.125	0.083
0.155	0.138	1	0.05	0.316
0.03	-0.125	0.05	1	0.239
0.121	0.083	0.316	0.239	1

Table 5

Rank	Student t-Test	Manual Selection	Correlation	Paired t-Test
1	SCF	SCF	SCF	SCF
2	MAD HOM	MAD HOM	MAD HOM	FGF-5
3	FGF-5	FGF-5	FGF-5	
4	FASLG	FASLG	FASLG	
5	MMP-10	MMP-10	PPY	
6	PPY	XPNPEP2	XPNPEP2	
7	XPNPEP2	FGF-21	FGF-21	
8	FGF-21	CXL17	MCP-3	
9	CXL17	MCP-3	FGF-BP1	
10	MCP-3	ESM-1	C15:0	
11	ESM-1	TNFB		
12	HK11	CTSV		
13	TRAIL	CD160		
14	FGF-BP1			
15	EN-RAGE			
16	C15:0			
17	TNFB			
18	CTSV			
19	ADA			
20	CD160			

Table 6

Rank	Random F	Manual S	Correlation	aired t-Test
1	SCF	SCF	SCF	SCF
2	MAD HOM	MAD HOM	MAD HOM	PPY
3	PPY	PPY	PPY	
4	FASLG	FASLG	FASLG	
5	FGF-5	FGF-5	FGF-5	
6	CXCL1	MMP-10	CXCL1	
7	MMP-10	XPNPEP2	XPNPEP2	
8	XPNPEP2	ESM-1	PHOSPHORIC ACID	
9	ESM-1	FLT3L	TLR3	
10	PHOSPHC	HK14	CD27	
11	PD-L1		FGF-BP1	
12	EPHA2			
13	FLT3L			
14	4E-BP1			
15	TRAIL			
16	MCP-1			
17	TLR3			
18	CD27			
19	FGF-BP1			
20	HK14			

Figure 1

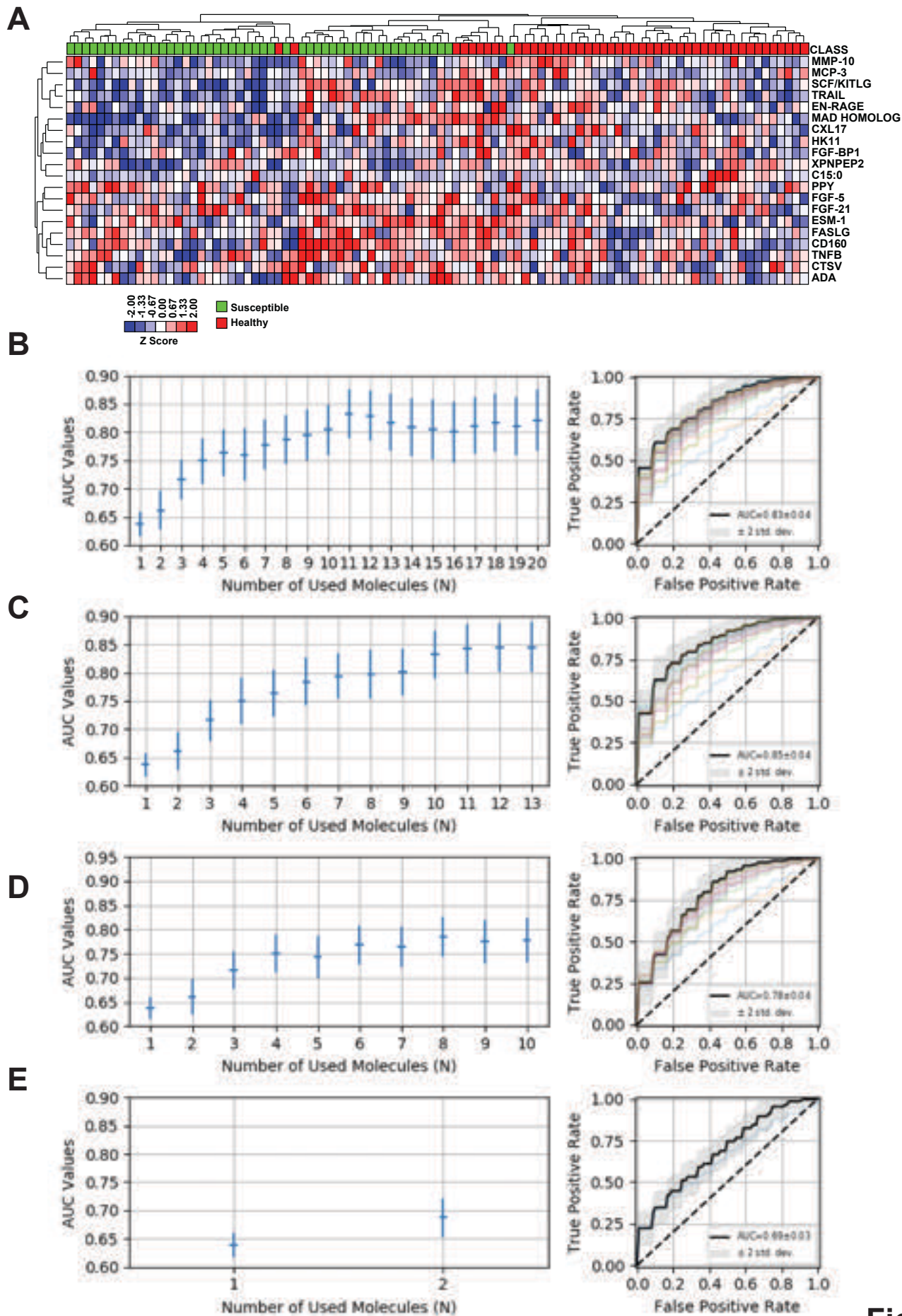


Figure 1

Figure 2

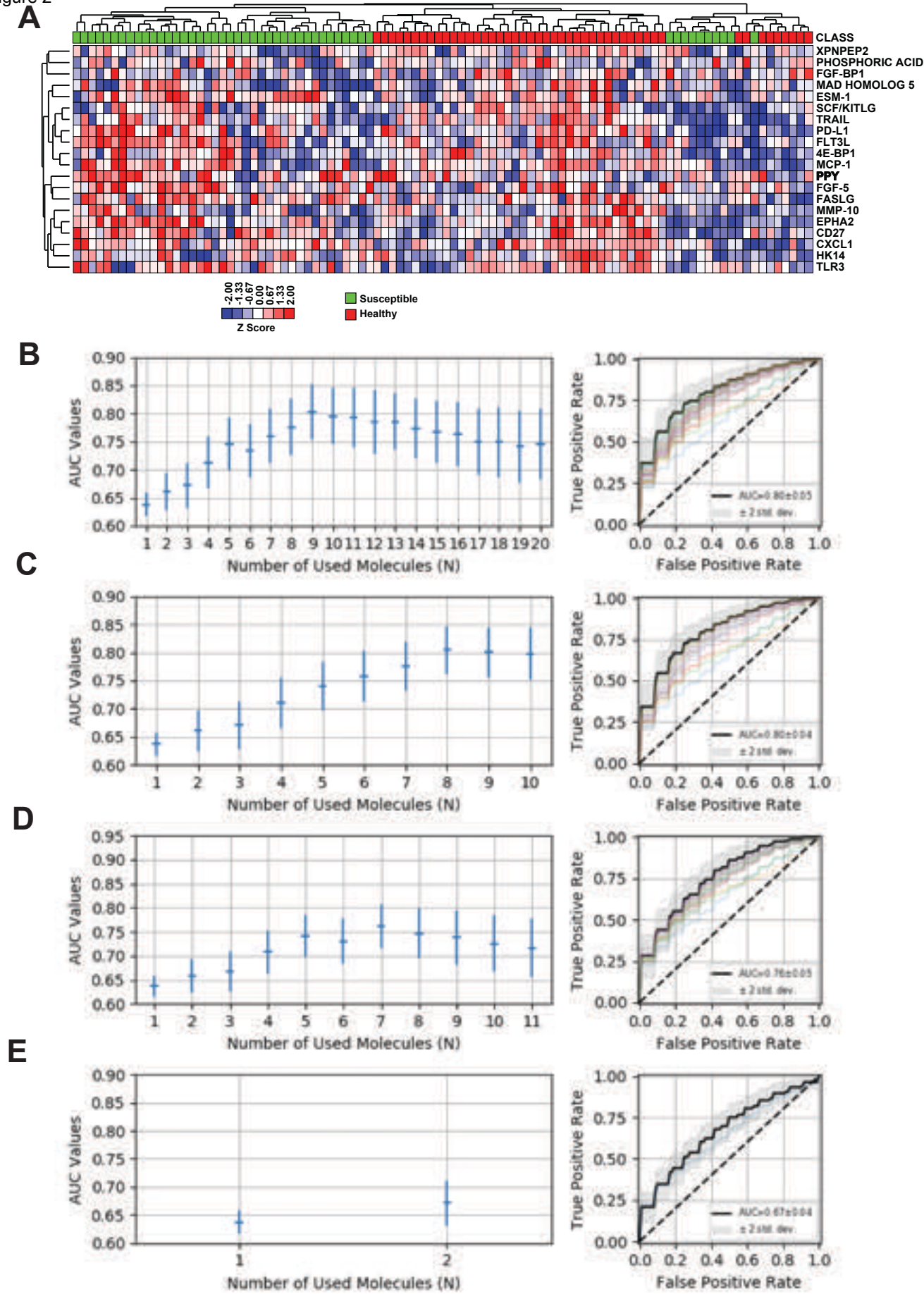


Figure 2

Figure 3

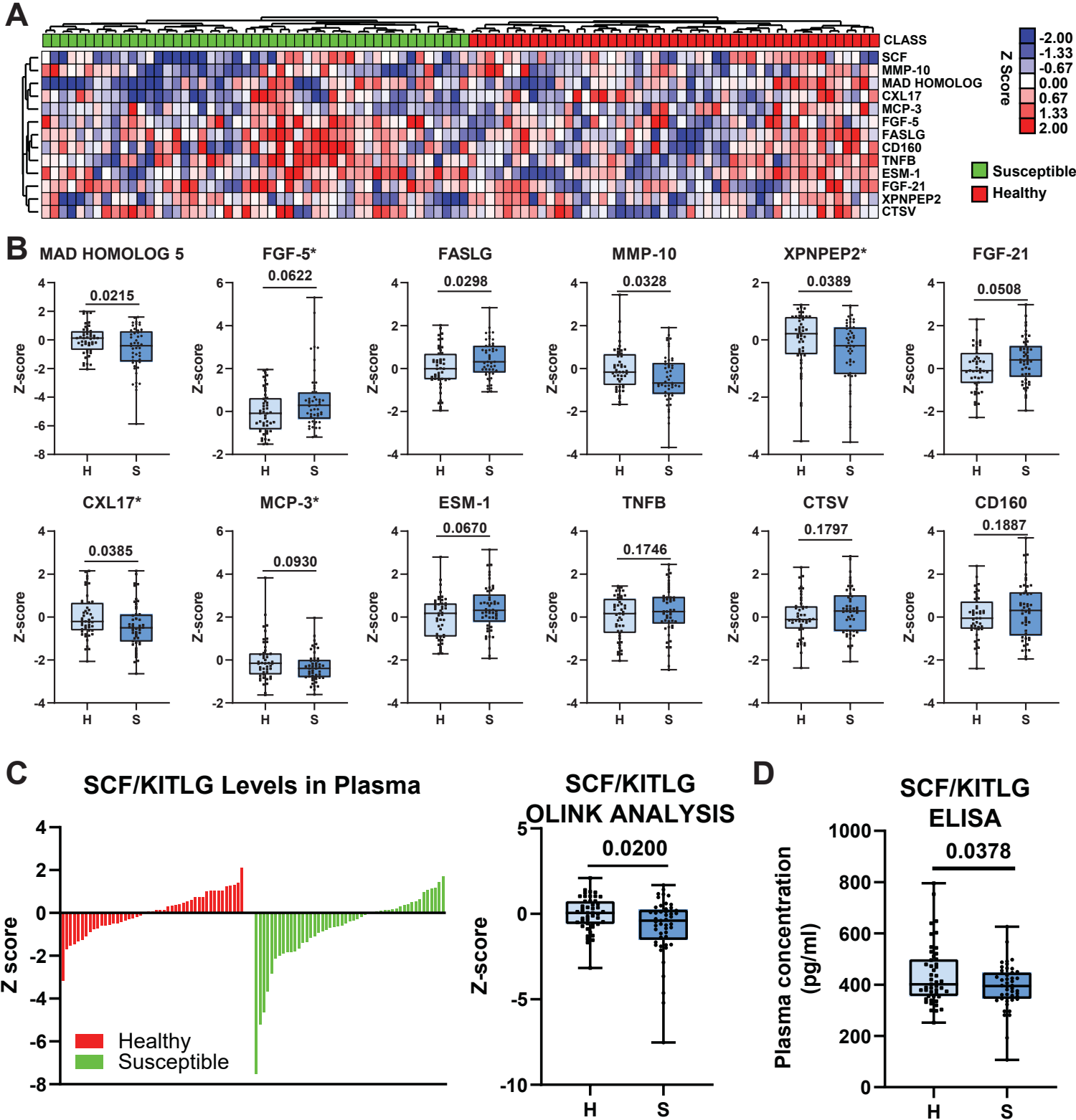


Figure 3